

When VLAD met Hilbert

Mehrtash Harandi

*NICTA and Australian National University
Canberra, Australia*

MEHRTASH.HARANDI@NICTA.COM.AU

Mathieu Salzmann

*NICTA and Australian National University
Canberra, Australia*

MATHIEU.SALZMANN@NICTA.COM.AU

Fatih Porikli

*NICTA and Australian National University
Canberra, Australia*

FATIH.PORIKLI@ANU.EDU.AU

Abstract

Vectors of Locally Aggregated Descriptors (VLAD) have emerged as powerful image/video representations that compete with or even outperform state-of-the-art approaches on many challenging visual recognition tasks. In this paper, we address two fundamental limitations of VLAD: its requirement for the local descriptors to have vector form and its restriction to linear classifiers due to its high-dimensionality. To this end, we introduce a kernelized version of VLAD. This not only lets us inherently exploit more sophisticated classification schemes, but also enables us to efficiently aggregate non-vector descriptors (*e.g.*, tensors) in the VLAD framework. Furthermore, we propose three approximate formulations that allow us to accelerate the coding process while still benefiting from the properties of kernel VLAD. Our experiments demonstrate the effectiveness of our approach at handling manifold-valued data, such as covariance descriptors, on several classification tasks. Our results also evidence the benefits of our nonlinear VLAD descriptors against the linear ones in Euclidean space using several standard benchmark datasets.

1. Introduction

This paper introduces several nonlinear formulations of *Vectors of Locally Aggregated Descriptors* (VLAD) that generalize their use to manifold-valued local descriptors, such as symmetric positive definite (SPD) matrices, and allows them to inherently exploit more sophisticated classification algorithms. Modern visual recognition techniques typically represent images by aggregating local descriptors, which, compared to image intensities, provide robustness to varying imaging conditions. From a historical point of view, this trend was gained momentum by the *Bag-of-Words* (BoW) model (Sivic et al., 2005; Grauman and Darrell, 2005; Lazebnik et al., 2006), which had a significant impact on recognition performance. Since then, the notable recent developments include dictionary-based solutions (Winn et al., 2005; Yang et al., 2009), Fisher Vectors (FV) (Perronnin and Dance, 2007; Perronnin et al., 2010b), VLAD (Jégou et al., 2010; Arandjelovic and Zisserman, 2013) and Convolutional Neural Networks (CNN) (Krizhevsky et al., 2012).

Among the aforementioned techniques, VLAD stands out for the following reasons:

- VLAD is computed via primitive operations. This makes VLAD extremely attractive when computational complexity is a concern.

- In contrast to CNNs, training a VLAD encoder is straightforward and not contingent on having a large training set.
- VLAD can be considered as a special case of FVs and hence inherits several properties of FVs. The most eminent one is its theoretical connection to the Fisher kernel (Jaakkola et al., 1999).
- From an empirical point of view, VLAD has been shown to either deliver state-of-the-art accuracy, or compete with the state-of-the-art methods. For instance, for scene classification on the MIT Indoor dataset, multi-scale VLAD, with only 4096 features, comfortably outperforms the mixture of FV and bag-of-parts, which relies on 221550 features (Gong et al., 2014).

Despite its unique features, VLAD comes with its own limitations. In particular, VLAD is designed to work with local descriptors in the form of vectors. Yet, several recent studies in computer vision suggest that structural data (*e.g.*, SPD matrices, graphs, orthogonal matrices) have the potential to provide more robust descriptors. Furthermore, since VLAD typically yields a high-dimensional image representation, it is mostly restricted to employing linear classifiers. Nonetheless, the effectiveness of kernel-based methods has been proven many a time in visual recognition (Gehler and Nowozin, 2009; Bo et al., 2010; Perronnin et al., 2010a; Vedaldi and Zisserman, 2012a).

In this paper, we present kernel based formulations of VLAD to address the aforementioned shortcomings. In particular, we first introduce a kernelized version of VLAD that relies on mapping of each local descriptor to a Reproducing Kernel Hilbert Space (RKHS). Since several valid kernel functions have recently been defined for non-vector data (Jayasumana et al., 2013; Harandi et al., 2014b), such a RKHS mapping can be applied to descriptors on different manifold topologies including SPD matrices and linear subspaces (Grassmannian). Having a RKHS mapping, we can aggregate VLAD over various geometries, thus, ultimately generalize the use of VLAD to local descriptors defined in non-vector spaces. Furthermore, the inherent nonlinearity of mappings to RKHS allows us to exploit more advanced classifiers with kernel VLAD.

In the spirit of computational efficiency, we also design three novel nonlinear approximations of our kernel VLAD; a Nyström method that obtains an explicit mapping to the Hilbert space, a local subspace-based representation of the data in Hilbert space, and a Fourier approximations based on the Bochner theorem. These approximations enjoy the similar properties of kernel VLAD, yet have the additional benefit of providing us with faster coding schemes. Interestingly, all our algorithms mostly preserve the simplicity of VLAD, in the sense that the extra computations merely consist of kernel evaluations potentially followed by projections (*i.e.*, matrix multiplications).

Table 1 provides a summary of the proposed methods and their attributes. Since each algorithm possesses unique features, it is not truly possible to pick one of the proposed methods as the ultimate winner. However, our experiments suggest that sVLAD is a good compromise between speed and accuracy, and would thus be our recommendation.

Our experimental evaluation demonstrate the effectiveness of our approach at handling manifold-valued data in a VLAD framework. Furthermore, we evidence the benefit of exploiting nonlinear classifiers for visual recognition by comparing the performance of our nonlinear VLAD with the standard one on several benchmark datasets, where the local descriptors have a vector form.

Table 1: Proposed methods and their properties. **Kernel** denotes the type of kernel function the algorithm can accept. For example, the *fVLAD* algorithm can only work with certain type of kernel functions while the *kVLAD* method can accept all type of kernel functions. **Coding** reflects the form of the output of the algorithm. For example, in the case of *kVLAD*, the output codes are only known implicitly. **Complexity** is the computational load.

Method	Kernel	Coding	Complexity
kVLAD	general	Implicit	High
nVLAD	general	Explicit	Low
fVLAD	specific	Explicit	Low
sVLAD	general	Explicit	Low

1.1 Related Work

Most of the popular image classification methods extract local descriptors (*i.e.*, at patch level), which are then aggregated into a global image representation (Lazebnik et al., 2006; Perronnin and Dance, 2007; Jégou et al., 2010; Perronnin et al., 2010b; Van Gemert et al., 2010; Krizhevsky et al., 2012; Arandjelovic and Zisserman, 2013).

When large amounts of training data are available, CNNs have now emerged as the method of choice to learn local descriptors. With limited number of training samples, existing methods typically opt for handcrafted features, such as SIFT.

To aggregate local features, in addition to operations such as average-pooling and max-pooling, histogram-based solutions (*e.g.*, BoW) have proven successful. Going beyond simple histograms has been an active topic of research for the past decade. For instance, (Lazebnik et al., 2006) aggregates histograms computed over different spatial regions. More recent developments, such as FVs (Perronnin and Dance, 2007) and VLAD (Jégou et al., 2010), suggest that high-order statistics should be encoded in the aggregation process.

In a separate line of research, structured descriptors (*e.g.*, covariance descriptors or linear subspaces) have been shown to provide robust visual models (Tuzel et al., 2008; Jayasumana et al., 2013; Harandi et al., 2013). Being of a non-vectorial form, aggregating such descriptors is hard to achieve beyond simple histograms. Nonetheless, one would like to benefit from the best of both worlds, that is, using robust non-vectorial descriptors in conjunction with state-of-the-art aggregation techniques, such as VLAD. This, in essence, is what we propose to achieve in this paper via kernelization. Furthermore, our approach has the additional advantage of allowing us to inherently exploit nonlinear classifiers that have proven powerful in visual recognition.

While a full review of kernel-based methods in computer vision is beyond the scope of this paper, the recent work of (Mairal et al., 2014) is of particular relevance here. (Mairal et al., 2014) introduces an approach to employing kernels within a CNN framework. Here, we perform a similar analysis within the VLAD framework, with the additional benefit of obtaining a representation that lets us work with manifold-valued data.

2. Nonlinear VLAD

In this section, we derive several nonlinear formulations of VLAD. To this end, we first start by reviewing the conventional VLAD and then discuss our approach to kernelizing it, followed by three approximations of the resulting kernel VLAD.

2.1 Conventional VLAD

Let $\mathcal{X} = \{\mathbf{x}_i\}_{i=1}^N$, $\mathbf{x}_i \in \mathbb{R}^d$ be a set of local descriptors extracted from a query image or a video. In VLAD (Jégou et al., 2010), the input space \mathbb{R}^d is partitioned into m Voronoi cells by means of a codebook \mathcal{C} with centers $\{\mathbf{c}_j\}_{j=1}^m$, $\mathbf{c}_j \in \mathbb{R}^d$. To obtain the codebook, the k-means algorithm is typically employed. Nevertheless, the use of supervised algorithms has also recently been advocated to build more discriminative codebooks (Peng et al., 2014). The VLAD code $\mathbf{v} \in \mathbb{R}^{md}$ for the query set \mathcal{X} is obtained by concatenating m Local Difference Vectors (LDV) δ_j storing, for each center, the sum of the differences between this center and each local descriptor assigned to this center. This can be written as

$$\mathbf{v}(\mathcal{X}) = \left[\delta_1^T(\mathcal{X}), \delta_2^T(\mathcal{X}), \dots, \delta_m^T(\mathcal{X}) \right]^T, \quad (1)$$

where

$$\delta_j(\mathcal{X}) = \sum_{i=1}^N a_j^i (\mathbf{c}_j - \mathbf{x}_i), \quad (2)$$

with a_j^i a binary weight encoding whether the local descriptor \mathbf{x}_i belongs to the Voronoi cell with center \mathbf{c}_j or not, *i.e.*, $a_j^i = 1$ if and only if the closest codeword to \mathbf{x}_i is \mathbf{c}_j .

2.2 Kernel VLAD (kVLAD)

As mentioned earlier, the conventional VLAD is designed to work with local descriptors of a vectorial form. As such, it cannot handle structured data representations, such as SPD matrices, or subspaces. While such representations could in principle be vectorized, this would (i) yield impractically high-dimensional VLAD vectors; and (ii) ignore the geometry of these structured representations, which has been demonstrated to result in accuracy losses (Pennec et al., 2006; Tuzel et al., 2006, 2008; Jayasumana et al., 2013). Here, we propose to address this problem by kernelizing VLAD.

To this end, let us redefine the query set of local descriptors as $\mathcal{X} = \{\mathbf{x}_i\}_{i=1}^N$, $\mathbf{x}_i \in \mathbb{X}$, where each descriptor lies in the space \mathbb{X} , which, in contrast to VLAD, is not restricted to be \mathbb{R}^d . In fact, the only constraint we impose is that \mathbb{X} comes with a valid positive definite *pd* kernel $k : \mathbb{X} \times \mathbb{X} \rightarrow \mathbb{R}$. For example, \mathbb{X} could be the space of SPD matrices, with the Gaussian kernel defined in (Sra, 2012; Jayasumana et al., 2013). According to the Moore-Aronszajn Theorem (Aronszajn, 1950), a *pd* kernel $k(\cdot, \cdot)$ induces a unique Hilbert space on \mathbb{X} , denoted hereafter by \mathcal{H} , with the property that there exists a mapping $\phi : \mathbb{X} \rightarrow \mathcal{H}$, such that $k(\mathbf{x}, \mathbf{y}) = \langle \phi(\mathbf{x}), \phi(\mathbf{y}) \rangle_{\mathcal{H}} = \phi(\mathbf{x})^T \phi(\mathbf{y})$. Here, we propose to make use of this property to map the local descriptors to \mathcal{H} , which is a vector space, and perform a VLAD-like aggregation in Hilbert space. The main difficulty arises from the fact that \mathcal{H} may be infinite-dimensional, and, more importantly, that the mapping ϕ corresponding to a given kernel k is typically unknown.

Let us suppose that we are given a codebook $\mathcal{C} = \{\phi(\mathbf{c}_i)\}_{i=1}^m$ in \mathcal{H} . For instance, this codebook can be computed using kernel kmeans. To compute a VLAD code in \mathcal{H} , we need to provide solutions for the following operations:

1. Determine the assignments $\{a_j^i\}$ in \mathcal{H} .
2. Express the LDVs in \mathcal{H} .

To determine the assignments, we note that

$$\|\phi(\mathbf{x}) - \phi(\mathbf{y})\|^2 = k(\mathbf{x}, \mathbf{x}) - 2k(\mathbf{x}, \mathbf{y}) + k(\mathbf{y}, \mathbf{y}) . \quad (3)$$

Therefore, for each local descriptor, the nearest codeword can be determined using kernel values only, *i.e.*, without having to know the mapping ϕ , which lets us directly define the assignments.

Unfortunately, expressing the LDVs in \mathcal{H} is not this straightforward. Clearly, the form of the LDVs, given by

$$\delta_j(\mathcal{X}) = \sum a_j^i \left(\phi(\mathbf{c}_j) - \phi(\mathbf{x}_i) \right) ,$$

with a_j^i obtained using Eq. 3, cannot be computed explicitly if the mapping ϕ is unknown, which is typically the case for popular kernels, such as RBF kernels. However, in most practical applications, the VLAD vector is not important by itself; What really matters for visual recognition is a notion of distance between two VLAD vectors. We therefore turn to the problem of computing the distance of two VLAD vectors in Hilbert space.

To this end, let $\mathcal{X} = \{\mathbf{x}_i\}_{i=1}^{N_{\mathcal{X}}}$, $\mathbf{x}_i \in \mathbb{X}$ and $\mathcal{Y} = \{\mathbf{y}_i\}_{i=1}^{N_{\mathcal{Y}}}$, $\mathbf{y}_i \in \mathbb{X}$ be two sets of local descriptors. The implicit VLAD code of \mathcal{X} in \mathcal{H} can be expressed as

$$\mathbf{v}_{\mathcal{H}}(\mathcal{X}) = \left[\delta_1^T(\mathcal{X}), \delta_2^T(\mathcal{X}), \dots, \delta_m^T(\mathcal{X}) \right]^T ,$$

and similarly for $\mathbf{v}_{\mathcal{H}}(\mathcal{Y})$. Now, we have

$$\begin{aligned} \langle \mathbf{v}_{\mathcal{H}}(\mathcal{X}), \mathbf{v}_{\mathcal{H}}(\mathcal{Y}) \rangle_{\mathcal{H}} &= \sum_{s=1}^m \delta_s^T(\mathcal{X}) \delta_s(\mathcal{Y}) \\ &= \sum_{s=1}^m \sum_{i=1}^{N_{\mathcal{X}}} \sum_{j=1}^{N_{\mathcal{Y}}} a_s^i a_s^j \left(\phi(\mathbf{c}_s) - \phi(\mathbf{x}_i) \right)^T \left(\phi(\mathbf{c}_s) - \phi(\mathbf{y}_j) \right) \\ &= \sum_{s=1}^m \sum_{i=1}^{N_{\mathcal{X}}} \sum_{j=1}^{N_{\mathcal{Y}}} a_s^i a_s^j \left(k(\mathbf{x}_i, \mathbf{y}_j) + k(\mathbf{c}_s, \mathbf{c}_s) - k(\mathbf{x}_i, \mathbf{c}_s) - k(\mathbf{y}_j, \mathbf{c}_s) \right) , \end{aligned} \quad (4)$$

which again only depends on kernel values.

With this inner product, a linear SVM, in its dual form, can directly be used for classification¹. In our experiments, we rely on this approach, which we refer to as kernel VLAD or **kVLAD** for short. This inner product, however, also allows us to employ an RBF-based kernel SVM, since

$$\|\mathbf{v}_{\mathcal{H}}(\mathcal{X}) - \mathbf{v}_{\mathcal{H}}(\mathcal{Y})\|^2 = \langle \mathbf{v}_{\mathcal{H}}(\mathcal{X}), \mathbf{v}_{\mathcal{H}}(\mathcal{X}) \rangle_{\mathcal{H}} - 2\langle \mathbf{v}_{\mathcal{H}}(\mathcal{X}), \mathbf{v}_{\mathcal{H}}(\mathcal{Y}) \rangle_{\mathcal{H}} + \langle \mathbf{v}_{\mathcal{H}}(\mathcal{Y}), \mathbf{v}_{\mathcal{H}}(\mathcal{Y}) \rangle_{\mathcal{H}} .$$

Note that this essentially yields two layers of kernels, *i.e.*, the RBF kernel of the SVM makes use of the distance, which itself is expressed in terms of kernel values.

While effective in practice, our kVLAD algorithm, as any kernel method, becomes computationally expensive when dealing with large datasets. In the remainder of this section, we therefore introduce three approximations to kVLAD, that address this limitation while still benefiting from the nice properties of kVLAD.

1. Note that this will yield a slightly different optimization problem than the standard kernel-based formulation, since in our case the inner product itself depends on several kernel values.

2.3 Nyström Approximation (nVLAD)

As a first approximation to kVLAD, we propose to make use of the Nyström method. Following (Perronnin et al., 2010a), this lets us obtain an explicit form for the mapping ϕ to the Hilbert space \mathcal{H} , and thus allows us to approximate a given kernel.

More specifically, let $\mathcal{T} = \{\mathbf{t}_i\}_{i=1}^M$, $\mathbf{t}_i \in \mathbb{X}$ be a collection of M training examples, and let \mathbf{K} be the corresponding kernel matrix, *i.e.*, $[\mathbf{K}]_{i,j} = k(\mathbf{t}_i, \mathbf{t}_j)$. We seek to approximate the elements of \mathbf{K} as inner products between r -dimensional vectors. In other words, we aim to find a matrix $\mathbf{Z} \in \mathbb{R}^{r \times M}$, such that $\mathbf{K} \simeq \mathbf{Z}^T \mathbf{Z}$. The best such approximation in the east-squares sense is given by $\mathbf{Z} = \mathbf{\Sigma}^{1/2} \mathbf{V}$, with $\mathbf{\Sigma}$ and \mathbf{V} the top r eigenvalues and corresponding eigenvectors of \mathbf{K} . From the Nyström method, for a new sample $\mathbf{x} \in \mathbb{X}$, the r -dimensional vector representation of the space induced by $k(\mathbf{x}, \cdot)$ can be written as

$$\mathbf{z}_N(\mathbf{x}) = \mathbf{\Sigma}^{-1/2} \mathbf{V} \left[k(\mathbf{x}, \mathbf{t}_1), \dots, k(\mathbf{x}, \mathbf{t}_M) \right]^T. \quad (5)$$

Given a set of local descriptors $\mathbb{X} = \{\mathbf{x}_i\}$, our **nVLAD** algorithm then consists of computing the corresponding $\{\mathbf{z}_N(\mathbf{x}_i)\}$, and making use of Eq. 1 and Eq. (2) with this new representation.

2.4 Local Subspace Approximation (sVLAD)

Here, we introduce a novel approximation of the Hilbert space \mathcal{H} based on the idea of local subspaces. To this end, we first note that the Nyström approximation yields one single global estimate of \mathcal{H} , used across all the codewords and all the descriptors. However, by looking at Eq. (2), we can see that the contribution of each codeword in the VLAD vector is independent of the other codewords, particularly since each local descriptor is assigned to a single codeword. Therefore, there is no reason for the approximation of \mathcal{H} to be shared across all the codewords and descriptors. This motivates us to define approximate spaces for each codeword individually.

To this end, let $\{\mathbf{t}_{s,j}\}_{j=1}^{N_s}$ be the set of training samples that generate the codeword \mathbf{c}_s . In other words, as in the conventional VLAD where $\mathbf{c}_s = \frac{1}{N_s} \sum_j \mathbf{t}_{s,j}$, we have $\phi(\mathbf{c}_s) = \frac{1}{N_s} \sum_i \phi(\mathbf{t}_{s,j})$. While, due to the unknown nature of ϕ , such a codeword cannot be explicitly computed, we can still evaluate the kernel function at this codeword, since

$$k(\mathbf{x}, \mathbf{c}_s) = \phi(\mathbf{x})^T \phi(\mathbf{c}_s) = \frac{1}{N_s} \sum_j \phi(\mathbf{x})^T \phi(\mathbf{t}_{s,j}) = \frac{1}{N_s} \sum_j k(\mathbf{x}, \mathbf{t}_{s,j}).$$

Here, we therefore propose to exploit the subspaces spanned by the training samples associated to each individual codeword to obtain an approximate representation of \mathcal{H} .

More specifically, let $\mathcal{S}_s = \text{span}(\{\phi(\mathbf{t}_{s,j})\}_{j=1}^{N_s})$. We then define

$$\bar{\delta}_s(\mathcal{X}) = \sum_{i=1}^N a_s^i \left(\pi_s(\phi(\mathbf{c}_s)) - \pi_s(\phi(\mathbf{x}_i)) \right), \quad (6)$$

with $\pi_s : \mathcal{H} \rightarrow \mathcal{S}_s$ the projection onto \mathcal{S}_s . These projections can be obtained following a similar intuition as for nVLAD. More precisely, let \mathbf{K}_s be the kernel matrix estimated from the training samples generating \mathbf{c}_s , *i.e.*, $[\mathbf{K}_s]_{i,j} = k(\mathbf{t}_{s,i}, \mathbf{t}_{s,j})$. By eigendecomposition, we can write $\mathbf{K}_s =$

$U_s \Lambda_s U_s^T$. Then, $\Phi_s U_s \Lambda_s^{-1/2}$, with $\Phi_s = [\phi(\mathbf{t}_{s,1}), \dots, \phi(\mathbf{t}_{s,N_s})]$, forms a basis for \mathcal{S}_s . As such, we can write

$$\pi_s(\mathbf{x}) = \Lambda_s^{-1/2} U_s \left[k(\mathbf{x}, \mathbf{t}_{s,1}), \dots, k(\mathbf{x}, \mathbf{t}_{s,N_s}) \right]. \quad (7)$$

The LDVs $\bar{\delta}_s(\mathcal{X})$ can then be obtained for all codeword \mathbf{c}_s , and concatenated to for the final **sVLAD** representation.

Remark 1 *Note that one can also use only the top r eigenvectors of \mathbf{K}_s to construct an r -dimensional local subspace in \mathcal{H} . This would not only yield the same dimensionality for all local subspaces, but could also potentially help discarding the noise associated to the $\{\mathbf{t}_{s,i}\}_{i=1}^{N_s}$.*

2.5 Fourier Approximation (fVLAD)

The previous two approximations apply to general kernels and both Euclidean and non-Euclidean data. In the Euclidean case, however, other approximations have been proposed for specific kernels (Rahimi and Recht, 2007; Vedaldi and Zisserman, 2012b). Since our experiments on Euclidean data all rely on RBF kernels, here, we discuss an approximation of this type of kernels based on the Bochner Theorem (Rudin, 2011).

According to the Bochner Theorem (Rudin, 2011), a shift-invariant kernel², such as Euclidean RBF kernel, can be obtained by the Fourier integral. As shown in (Rahimi and Recht, 2007), for real-valued kernels, this can be expressed as

$$k(\mathbf{x}_i - \mathbf{x}_j) = \int_{\mathbb{R}^d} p(\omega) z_F(\mathbf{x}_i) z_F(\mathbf{x}_j) d\omega, \quad (8)$$

where $z_F(\mathbf{x}) = \sqrt{2} \cos(\omega^T \mathbf{x} + b)$, with b a random variable drawn from $[0, 2\pi]$. In other words, $k(\mathbf{x}_i, \mathbf{x}_j) = k(\mathbf{x}_i - \mathbf{x}_j)$ is the expected value of $z_F(\mathbf{x}_i) z_F(\mathbf{x}_j)$ under the distribution $p(\omega)$. For the RBF kernel $k(\mathbf{x}_i, \mathbf{x}_j) = \exp(-\|\mathbf{x}_i - \mathbf{x}_j\|^2 / 2\sigma^2)$, we have $p(\omega) = \mathcal{N}(0, \sigma^{-2} \mathbf{I}_d)$.

Let $\{\omega_i\}_{i=1}^r$, $\omega_i \in \mathbb{R}^d$, be i.i.d. samples drawn from the normal distribution $\mathcal{N}(0, \sigma^{-2} \mathbf{I}_d)$, and $\{b_i\}_{i=1}^r$ be samples uniformly drawn from $[0, 2\pi]$. Then, the r dimensional estimate of $\phi(\mathbf{x}) \in \mathcal{H}$ is given by

$$z_F(\mathbf{x}) = \sqrt{\frac{2}{r}} \left[\cos(\omega_1^T \mathbf{x} + b_1), \dots, \cos(\omega_r^T \mathbf{x} + b_r) \right]. \quad (9)$$

Similarly to nVLAD, we can then compute $z_F(\mathbf{x}_i)$ for each local descriptor \mathbf{x}_i , and use Eq. (1) and Eq. (2) to obtain a code. In our experiments, we refer to this approach, which only applies to Euclidean data, as **fVLAD**.

2.6 Further Considerations

Normalization:

Recent developments have suggested that the discriminatory power of VLAD could be boosted by additional post-processing steps, such as ℓ_2 power normalization and signed square rooting normalization (Arandjelovic and Zisserman, 2013; Gong et al., 2014). The ℓ_2 power normalization, where each block in VLAD is normalized individually, can easily be performed in kVLAD, since

$$\|\delta_s(\mathbb{X})\|_{\mathcal{H}}^2 = \sum_{i,j=1}^{N_s} a_s^i a_s^j \left(k(\mathbf{x}_i, \mathbf{x}_j) + k(\mathbf{c}_s, \mathbf{c}_s) - k(\mathbf{x}_i, \mathbf{c}_s) - k(\mathbf{x}_j, \mathbf{c}_s) \right)$$

2. A kernel function is shift invariant if $k(\mathbf{x}_i, \mathbf{x}_j) = k(\mathbf{x}_i - \mathbf{x}_j)$.

is only dependent on kernel values. As a result, the inner product of Eq. 4 after normalizing each VLAD block independently, *i.e.*,

$$\left\langle \bar{\mathbf{v}}_{\mathcal{H}}(\mathbb{X}), \bar{\mathbf{v}}_{\mathcal{H}}(\mathbb{Y}) \right\rangle_{\mathcal{H}} = \sum_{s=1}^k \frac{\left\langle \delta_s(\mathbb{X}), \delta_s(\mathbb{Y}) \right\rangle}{\|\delta_s(\mathbb{X})\|_{\mathcal{H}} \|\delta_s(\mathbb{Y})\|_{\mathcal{H}}} ,$$

will also only depend on kernel values. By contrast, however, the signed square rooting normalization can only be achieved when explicit forms of the descriptors are available, *i.e.*, in nVLAD, sVLAD and fVLAD.

Kernelizing Fisher Vectors:

Due to the connection between VLAD and FVs, it seems natural to rely on the ideas discussed above to kernelize FVs. One difficulty in kernelizing FV, however, arises from the fact that Gaussian distributions, which are required to model the probability distributions in FV, are not well-defined in RKHS. More specifically, to fit a Gaussian distribution in a d -dimensional space, at least d independent observations (training samples) are required, to ensure that the covariance matrix of the distribution is not rank deficient. Obviously, for an infinite dimensional RKHS, this requirement cannot be met. While, in principle, it is possible to regularize the distributions, *e.g.*, (Zhou and Chellappa, 2006), we believe that an in-depth analysis of this approach to kernelize FVs goes beyond the scope of this paper. Note, however, that our approximations of \mathcal{H} can be applied verbatim to derive approximate formulations of kernel FV.

3. Experiments

We now evaluate our different algorithms, *i.e.*, kVLAD, nVLAD, sVLAD and fVLAD, on several recognition tasks. As mentioned before, our main motivation for this work was to be able to exploit the power of the VLAD aggregation scheme to tackle problems where the input data is not in vectorial form. Therefore, we focus on two such types of data, which have become increasingly popular in computer vision, namely Covariance Descriptors (CovDs), which lie on SPD manifolds, and linear subspaces which form Grassmann manifolds. Nevertheless, in addition to this manifold-valued data, we also evaluate our algorithms in Euclidean space.

3.1 SPD Manifold

In computer vision, SPD matrices have been shown to provide powerful representations for images and videos via region covariances (Tuzel et al., 2006). Such representations have been successfully employed to categorize, *e.g.*, textures (Tuzel et al., 2006; Harandi et al., 2014a), pedestrians (Tuzel et al., 2008) and faces (Harandi et al., 2014a).

SPD matrices can be thought of as an extension of positive numbers and form the interior of the positive semidefinite cone. It is possible to directly employ the Frobenius norm as a similarity measure between SPD matrices, hence analyzing problems involving such matrices via Euclidean geometry. However, as several studies have shown, undesirable phenomena may occur when Euclidean geometry is utilized to manipulate SPD matrices (Pennec et al., 2006; Tuzel et al., 2008; Jayasumana et al., 2013). Here, instead, we make use of the Stein divergence defined as

$$\delta_S^2(\mathbf{A}, \mathbf{B}) = \ln \det \left(\frac{\mathbf{A} + \mathbf{B}}{2} \right) - \frac{1}{2} \ln \det (\mathbf{A}\mathbf{B}) . \quad (10)$$

This divergence was shown to yield a positive definite Gaussian kernel (Sra, 2012), named the Stein kernel given by $k_S : \mathcal{S}_{++}^n \times \mathcal{S}_{++}^n \rightarrow \mathbb{R}$ such that $k_S(\mathbf{A}, \mathbf{B}) = \exp(-\sigma \delta_S^2(\mathbf{A}, \mathbf{B}))$. In all our experiments on SPD manifolds, the bandwidth of this kernel was determined by cross-validation on the training data.

A standard approach when dealing with an SPD manifold consists of flattening the manifold using the diffeomorphism $\log : \mathcal{S}_{++}^n \rightarrow \text{Sym}(n)$, where \log and $\text{Sym}(n)$ denote the principal matrix logarithm and the space of symmetric matrices of size n , respectively. Given that $\text{Sym}(n)$ is a vector space, one can then directly employ tools from Euclidean geometry, here the VLAD algorithm, to analyze SPD matrices mapped to that space. In our experiments, we refer to this baseline as log-Euclidean VLAD or *IE-VLAD* following the terminology used in (Arsigny et al., 2006). Note that this strategy has been successfully employed in several recent studies (e.g., for image semantic segmentation (Carreira et al., 2012)).

Furthermore, we also compare our algorithms against the state-of-the-art Weighted ARray of COvariances (WARCO) (Tosato et al., 2013), Covariance Discriminative Learning (CDL) (Wang et al., 2012) and Riemannian Sparse Representation using the Stein divergence (RSR-S) (Harandi et al., 2015) algorithms. In WARCO, an image is decomposed into a number of overlapped patches, each of which is represented with a CovD. Classification is then performed by combining the output of a set of kernel classifiers trained on local patches. In essence, WARCO pursues the same goal as us, *i.e.*, to aggregate local non-vectorial descriptors, which makes it probably the most relevant baseline, here. By contrast, following (Wang et al., 2012; Harandi et al., 2015), we have used both CDL and RSR-S holistically, *i.e.*, every image was described by one SPD matrix.

In the following experiments on the SPD manifold, we used a codebook of size 32 for all variants of the VLAD algorithm. Empirically, we observed that, for any algorithm, larger codebooks did not significantly improve the performance. To provide a fair comparison against WARCO, we use the same set of features as (Tosato et al., 2013). More specifically, from a local patch, a 13×13 CovD is extracted using the features

$$f(x, y) = [h_1(Y), \dots, h_8(Y), Y, C_b, C_r, \|g(Y)\|, \angle(g(Y))]^T,$$

where $f(x, y)$ denotes the feature vector at location (x, y) and Y , C_b and C_r are the three color channels from the CIE Lab color space at (x, y) . $h_i(\cdot)$ is the scaled symmetric Difference Of Offset Gaussian filter bank, and $\|g(Y)\|$ and $\angle(g(Y))$ are the gradient magnitude and orientation calculated on the Y channel (see (Tosato et al., 2013) for details). The same set of features was used for CDL and RSR-S.

Head Orientation Classification. As a first experiment, we consider the problem of classifying head orientation using the *QMUL* and *HOCoffee* datasets (Tosato et al., 2013). The *QMUL* head dataset contains 19292 images of size 50×50 , captured in an airport terminal. The *HOCoffee* dataset (see Fig. 1 for examples) contains 18117 head images of size 50×50 . The images typically include a margin of 10 pixels on average, so that the actual average dimension of the heads is 30×30 pixels. Both datasets come with predefined training and test samples.

In Table. 2, we report the performance of kVLAD, sVLAD and nVLAD, as well as of WARCO and IE-VLAD, on the *QMUL* and *HOCoffee* datasets. Note that kVLAD and sVLAD both yield higher accuracies than the state-of-the-art WARCO algorithm. For example, on *HOCoffee*, the accuracy of kVLAD surpasses that of WARCO by more than 5%. Note also that kVLAD and sVLAD yield very similar accuracies, which evidences the good quality of our local subspace approximation. Interestingly, sVLAD even outperforms kVLAD on *QMUL*. This can be attributed to the



Figure 1: Samples from the HOCoffee dataset.



Figure 2: Samples from the HOC dataset.

square root normalization, which is not possible for kVLAD. Without this normalization, the performance of sVLAD drops by roughly 1%, and thus remains close to, but slightly lower than that of kVLAD. Among the approximations, sVLAD is superior to nVLAD. This is not really surprising, since nVLAD relies on a single subspace for all its codewords, whereas sVLAD exploits more local representations.

Body Orientation Classification. As a second task on the SPD manifold, we consider the problem of determining body orientation from images using the Human Orientation Classification (HOC) dataset (Tosato et al., 2013). The HOC dataset contains 11881 images of size 64×32 (see Fig. 2 for examples) and comprises 4 orientation classes (Front, Back, Left, and Right). In Table. 2, we compare the performance of kVLAD, sVLAD and nVLAD against that of WARCO and IE-VLAD. First, we note that all VLAD variants, including IE-VLAD, are superior to the WARCO algorithm. This demonstrates the effectiveness of the VLAD aggregation scheme. Moreover, we note that all our algorithms outperform IE-VLAD. The highest accuracy is obtained by sVLAD which again, in comparison to kVLAD, benefits from the square root normalization.

Altogether, our experiments on SPD manifolds demonstrate that our approach offers an attractive solution to exploiting the information from local patches. Note that, except for a handful of studies (*e.g.*, WARCO), CovDs are usually extracted from entire images, hence making them questionable for challenging classification tasks. This is typically due to the fact that aggregating non-vectorial is an open problem, to which we provide a solution in this paper.

3.2 Grassmann Manifold

The space of p dimensional subspaces in \mathbb{R}^d for $0 < p \leq d$ is not a Euclidean space, but a Riemannian manifold known as the Grassmann manifold $\mathcal{G}(p, d)$. A point $\mathcal{U} \in \mathcal{G}(p, d)$ is typically represented by a $d \times p$ matrix \mathbf{U} with orthonormal columns, such that $\mathcal{U} = \text{Span}(\mathbf{U})$. The choice of the basis to represent \mathcal{U} is arbitrary and metrics on $\mathcal{G}(p, d)$ are defined so as to be invariant to this

Method	QMUL	HOCoffee	HOC
WARCO (Tosato et al., 2013)	91%	80%	78%
CDL (Wang et al., 2012)	81.6%	71.2%	77.8%
RSR-S (Harandi et al., 2015)	82.7%	65.9%	78.9%
IE-VLAD	87.6%	77.2%	79.7%
nVLAD	88.9%	83.4%	81.4%
sVLAD	92.7%	84.0%	84.1%
kVLAD	92.2%	85.3%	83.1%

Table 2: Recognition accuracies for QMUL, HOCoffe and HOC.

choice. The projection distance is a typical choice of such metric. It was recently shown to induce a valid positive definite kernel on $\mathcal{G}(p, d)$ (Harandi et al., 2014b), *i.e.*, the projection RBF kernel defined as

$$k_p(\mathbf{A}, \mathbf{B}) = \exp(\sigma \|\mathbf{A}^T \mathbf{B}\|_F^2), \sigma > 0. \quad (11)$$

As for the SPD manifold, in our experiments, the bandwidth of this kernel was obtained by cross-validation on the training data.

Several state-of-the-art image-set matching methods model sets of images as subspaces (Harandi et al., 2013, 2014b). However, to the best of our knowledge, all these methods rely on a holistic subspace representation. This again is probably due to the fact that, before this paper, no aggregation schemes on Grassmann manifolds have ever been proposed. Our approach, by contrast, enables us to break an image-set into smaller blocks, represent each block by a linear subspace, and aggregating these subspace to form a complete image-set descriptor.

In our experiments, we compare the results of our algorithms against four baselines: First, similarly to the log-Euclidean approach on SPD manifolds, we propose to flatten $\mathcal{G}(p, d)$ at $\mathbf{I}_{d \times p}$ ³ and perform conventional VLAD in the resulting Euclidean space. We refer to this method as IE-VLAD. As a second baseline, we make use of the state-of-the-art Grassmannian Sparse Coding (gSC) algorithm of (Harandi et al., 2013), which describes each image-set with a single linear subspace. We also employ the kernel version of the Affine Hull Method (kAHM) introduced in (Cevikalp and Triggs, 2010) and the CDL algorithm (Wang et al., 2012) as other state-of-the-art baselines for image-set matching. Below, we evaluate the performance of our algorithms and of the baselines on three different classification problems, *i.e.*, object recognition, action classification and pose categorization from image-sets.

Action Recognition. As a first experiment on the Grassmannian, we make use of the Ballet dataset (Wang and Mori, 2009). The Ballet dataset consists of 8 complex motion patterns performed by 3 subjects (see Fig. 3 for examples). We extracted 1200 image-sets by grouping 5 frames depicting the same action into one image-set. The local descriptors for each image-set were obtained by splitting the set into small blocks of size $32 \times 32 \times 3$ and utilizing Histogram of Oriented Gradient (HOG) (Dalal and Triggs, 2005). We then created subspaces of size 31×3 , hence points on $\mathcal{G}(3, 31)$. We randomly chose 50% of imagesets for training and used the remaining sets as test samples. The process of random splitting was repeated ten times and the average classification accuracy is reported.

3. We use $\mathbf{I}_{d \times p}$ to denote the truncated identity matrix.



Figure 3: Samples from the Ballet dataset.(grayscale images were used in our experiments).

In Table 3, we report the accuracy of algorithms and of the gSC and IE-VLAD baselines. First, note that all the local approaches outperform the holistic gSC method. Furthermore, similarly to the two experiments on SPD manifolds, the maximum accuracy is obtained by sVLAD, closely followed by kVLAD.

Given the simplicity of the IE-VLAD method, it is interesting to verify if it can measure up to our kernel extensions by enlarging its dictionary. To this end, we increased the size of the dictionary in IE-VLAD up to the point where the performance started to decrease (256 atoms). While this indeed improved the accuracy of IE-VLAD up to the best accuracy of 91.7%, it remains significantly below the performance of sVLAD.

Object Recognition. For the task of object recognition from image-sets, we used the CIFAR dataset (Krizhevsky and Hinton, 2009). The CIFAR dataset contains 60000 images (32×32 pixels) from 10 different object categories. From this dataset, we generated 6000 image-sets, each one containing 10 random images of the same object. In our experiments, we used 1500 image-sets for training and the remaining 4500 image-sets as test data. We report accuracies averaged over 10 random image-set generation processes.

To generate local descriptors, we decomposed each image-set into small blocks of size $8 \times 8 \times 5$. Each block was then represented by a point on $\mathcal{G}(5, 64)$ using SVD. In Table. 3, we compare the results of kVLAD, sVLAD and nVLAD against those of IE-VLAD and gSC. Here, kVLAD yields the best accuracy followed by sVLAD.

Pose Classification. As a last experiment on the Grassmannian, we evaluated the performance of our algorithms on the task of pose categorization using the CMU-PIE face dataset (Sim et al., 2003). The CMU-PIE face dataset contains images of 67 subjects under 13 different poses and 21 different illuminations (see Fig. 4 for examples). The images were closely cropped to enclose the face region and resized to 64×64 . We extracted 1700 image-sets by grouping 6 images with the same pose, but different illuminations into one image-set. The local descriptors for each image set were obtained by splitting the set into small blocks of size $32 \times 32 \times 3$ from which we computed Histogram of LBP (Ojala et al., 2002). We then created subspaces of size 58×3 , hence points on $\mathcal{G}(3, 58)$. Table 3 compares the results of nVLAD, sVLAD and kVLAD against those of gSC and IE-VLAD. The highest accuracy is obtained by kVLAD, this time by a large margin over the second best, sVLAD. Note that, with this dataset, flattening the manifold through its tangent space at $\mathbf{I}_{58 \times 3}$ seems to incur strong distortions, as indicated by low performance of IE-VLAD.

Remark 2 Several recent studies (e.g., (Wang et al., 2012; Huang et al., 2015)) have tackled the problem of image-set matching using the geometry of SPD manifolds via covariance descriptors. Table 3, however, suggests that, for our experiments, the resulting global covariance descriptors do



Figure 4: Samples from CMU-PIE.

Method	Ballet	CIFAR	CMU-PIE
gSC (Harandi et al., 2013)	79.7%	59.9%	75.5%
kAHM (Cevikalp and Triggs, 2010)	85.8%	36.1%	55.3%
CDL (Wang et al., 2012)	73.1%	54.7%	64.6%
IE-VLAD	91.1%	46.2%	59.6%
nVLAD	88.9%	62.2%	79.5%
sVLAD	94.4%	65.2%	80.1%
kVLAD	92.2%	67.9%	86.3%

Table 3: Accuracies for Ballet, CIFAR and CMU-PIE.

not measure up to subspaces, as evidenced by the performance of CDL in comparison to gSC. We conjecture that this is due to the small number of images in each set, which makes the SPD matrices rank deficient (regularization was used to overcome this issue) and less discriminative. Interestingly, however, we also evaluated sVLAD using local SPD matrices instead of local subspace, and achieved an accuracy of 93.4% on the Ballet dataset. While this remains slightly below what Grassmannian geometry can achieve, it clearly shows the strength of our framework, which, by using a local representation, outperforms the global descriptors of CDL by more than 20%.

3.3 Euclidean Space

Our final experiments are devoted to Euclidean spaces. To this end, we compare the performance of sVLAD, fVLAD and nVLAD against the conventional VLAD (implementation provided in (Vedaldi and Fulkerson, 2008)) on Pascal VOC 2007 (Everingham et al., 2010) and on the Flickr Material

Method	mAP
SPM (Lazebnik et al., 2006)	54.3%
OCP (Russakovsky et al., 2012)	57.2%
Sup-VLAD (Peng et al., 2014)	60.9%
VLAD	54.7%
nVLAD	56.2%
fVLAD	55.8%
sVLAD	60.3%

Table 4: mean Average Precision (mAP) for VOC 2007 dataset.

Method	CCR
aLDA (Sharan et al., 2013)	44.6%
MS4C (Li, 2014)	50.0%
DTD_{RBF} (Cimpoi et al., 2014)	53.1%
VLAD	49.4%
nVLAD	52.3%
fVLAD	50.3%
sVLAD	55.2%

Table 5: Correct Classification Rate (CCR) for FMD dataset.



Figure 5: Examples of the FMD texture dataset.

Database (FMD) (Sharan et al., 2013) (see Fig. 5 for examples) . Pascal VOC 2007 (Everingham et al., 2010) contains 9963 images from 20 object categories. The FMD contains 1000 images from 10 different material categories (Sharan et al., 2013). Both datasets have been extensively used to benchmark coding techniques.

In our experiments, we realized that the computational load of kVLAD becomes overwhelming on Pascal VOC07 and FMD as a result of large amount of local descriptors. Hence, we will only report the performance of nVLAD, fVLAD and sVLAD here. The size of codebooks was set to 256 and SIFT descriptors (with whitening) were considered as local features. For fVLAD and nVLAD, the size of the RKHS was chosen to be 256 (almost 3 times larger than the original space). While increasing the dimensionality of the RKHS could potentially improve the results, it would come at the expense of increasing the computational burden of coding.

Table 4 compares the recognition accuracies of the proposed coding techniques with conventional VLAD, Spatial Pyramid Matching (SPM) (Lazebnik et al., 2006), Object-Centric spatial Pooling (OCP) (Russakovsky et al., 2012) and supervised dictionary learning for VLAD (Sup-VLAD) (Peng et al., 2014). Similarly to our experiments on manifolds, sVLAD outperforms the fixed approximation techniques (*i.e.*, fVLAD and nVLAD). Importantly, we observe that our three algorithms outperform traditional methods such as SPM and VLAD. Furthermore, sVLAD also outperforms the state-of-the-art pooling method OCP (Russakovsky et al., 2012), and performs on par with the supervised Sup-VLAD. This latter comparison motivates an interesting future research direction to learn a supervised dictionary in RKHS.

Table 5 compares the recognition accuracies of nVLAD, fVLAD and sVLAD against VLAD and the state-of-the-art methods augmented Latent Dirichlet Allocation (aLDA) (Sharan et al., 2013), Multi-Scale Spike-and-Slab Sparse Coding (MS4C) (Li, 2014), and Describable attributes (DTD_{RBF}) (Cimpoi et al., 2014) on the FMD dataset. In essence, we can see that (i) our algorithms outperform VLAD, with sVLAD the best-performing method; (ii) our algorithms outperform the state-of-the-art aLDA and MS4C methods; and (iii) while DTD yields higher accuracy than our fixed approximations (*i.e.*, nVLAD and fVLAD), it is still outperformed by our sVLAD algorithm.

3.4 Encoding times

Before concluding, we provide the coding times for the proposed methods on the three different geometries studied in this work. In particular, we measured the encoding times of sVLAD, fVLAD and nVLAD on a Quad-core machine using Matlab. We also measured the running time to compute Eq. 4, which shows the computational load of kVLAD.

The parameters values of the algorithms when measuring these timings were those used in our experiments. More specifically, for the SPD and Grassmann manifolds, the size of codebook was chosen to be 32, while, in the case of Euclidean space, it was set to 256. Note that for the Euclidean

Method	SPD	Grassmann	Euclidean
nVLAD	650ms	1600ms	35ms
fVLAD	N/A	N/A	100ms
sVLAD	750ms	1700ms	950ms
kVLAD	80ms	155ms	45ms

Table 6: Running times for fVLAD, nVLAD, sVLAD and kVLAD on three different geometries. Note that the running times for fVLAD, nVLAD and sVLAD show the coding time for an image/video, while, in the case of kVLAD where not explicit encoding is performed, it shows the time needed to evaluate Eq. 4

case, we assumed that 1000 local descriptors were computed on each image, while, for the SPD and Grassmann manifolds, this number was set to 100. Table 6 reports all the running times.

4. Conclusions and Future Work

In this paper, we have introduced a kernel extension of the VLAD encoding scheme. We have also proposed several approximations to this kernel formulation in the interest of speeding up the encoding process. Not only do the resulting algorithm let us exploit more sophisticated classification schemes in the VLAD framework, but they also allow us to aggregate local descriptors that do not lie in Euclidean space. Our experiments have evidenced that our algorithms outperform state-of-the-art methods, such as WARCO (Tosato et al., 2013) and gSC (Harandi et al., 2013), on several manifold-based recognition tasks. Furthermore, they have also shown that our new encoding schemes yield superior results compared to the conventional VLAD algorithm. In the future, we plan to explore possible ways of kernelizing the Fisher vector (Perronnin and Dance, 2007) method. We also intend to study the concept of coresets (Har-Peled and Mazumdar, 2004) to reduce the computational complexity of coding.

References

- Relja Arandjelovic and Andrew Zisserman. All about vlad. In *CVPR*, 2013.
- Nachman Aronszajn. Theory of reproducing kernels. *Transactions of the American mathematical society*, 1950.
- Vincent Arsigny, Olivier Commowick, Xavier Pennec, and Nicholas Ayache. A log-euclidean framework for statistics on diffeomorphisms. In *MICCAI*. 2006.
- Liefeng Bo, Xiaofeng Ren, and Dieter Fox. Kernel descriptors for visual recognition. In *NIPS*, 2010.
- Joao Carreira, Rui Caseiro, Jorge Batista, and Cristian Sminchisescu. Semantic segmentation with second-order pooling. In *ECCV*. 2012.
- Hakan Cevikalp and Bill Triggs. Face recognition based on image sets. In *CVPR*, pages 2567–2573. IEEE, 2010.

- Mircea Cimpoi, Subhrajyoti Maji, Iasonas Kokkinos, Salina Mohamed, and Andrea Vedaldi. Describing textures in the wild. In *Computer Vision and Pattern Recognition (CVPR), 2014 IEEE Conference on*, pages 3606–3613, 2014.
- Navneet Dalal and Bill Triggs. Histograms of oriented gradients for human detection. In *CVPR*, 2005.
- Mark Everingham, Luc Van Gool, Christopher KI Williams, John Winn, and Andrew Zisserman. The pascal visual object classes (voc) challenge. *IJCV*, 88(2), 2010.
- Peter Gehler and Sebastian Nowozin. On feature combination for multiclass object classification. In *CVPR*, 2009.
- Yunchao Gong, Liwei Wang, Ruiqi Guo, and Svetlana Lazebnik. Multi-scale orderless pooling of deep convolutional activation features. In *ECCV*. 2014.
- Kristen Grauman and Trevor Darrell. The pyramid match kernel: Discriminative classification with sets of image features. In *ICCV*, 2005.
- Sariel Har-Peled and Soham Mazumdar. On coresets for k-means and k-median clustering. In *ACM symposium on Theory of computing*, 2004.
- Mehrtash Harandi, Conrad Sanderson, Chunhua Shen, and Brian Lovell. Dictionary learning and sparse coding on grassmann manifolds: An extrinsic solution. In *ICCV*, 2013.
- MehrtashT. Harandi, Mathieu Salzmann, and Richard Hartley. From manifold to manifold: Geometry-aware dimensionality reduction for spd matrices. In *ECCV*. 2014a.
- MehrtashT. Harandi, Mathieu Salzmann, Sadeep Jayasumana, Richard Hartley, and Hongdong Li. Expanding the family of grassmannian kernels: An embedding perspective. In *ECCV*. 2014b.
- M.T. Harandi, R. Hartley, B. Lovell, and C. Sanderson. Sparse coding on symmetric positive definite manifolds using bregman divergences. *TNNLS*, PP(99):1–1, 2015. ISSN 2162-237X. doi: 10.1109/TNNLS.2014.2387383.
- Zhiwu Huang, Ruiping Wang, Shiguang Shan, and Xilin Chen. Face recognition on large-scale video in the wild with hybrid euclidean-and-riemannian metric learning. *Pattern Recognition*, 48(10):3113 – 3124, 2015. ISSN 0031-3203. Discriminative Feature Learning from Big Data for Visual Recognition.
- Tommi Jaakkola, David Haussler, et al. Exploiting generative models in discriminative classifiers. In *NIPS*, 1999.
- Sadeep Jayasumana, Richard Hartley, Mathieu Salzmann, Hongdong Li, and Mehrtaash Harandi. Kernel methods on the riemannian manifold of symmetric positive definite matrices. In *CVPR*, 2013.
- Hervé Jégou, Matthijs Douze, Cordelia Schmid, and Patrick Pérez. Aggregating local descriptors into a compact image representation. In *CVPR*, 2010.

- Alex Krizhevsky and Geoffrey Hinton. Learning multiple layers of features from tiny images. *Tech. Rep*, 2009.
- Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In *NIPS*, 2012.
- Svetlana Lazebnik, Cordelia Schmid, and Jean Ponce. Beyond bags of features: Spatial pyramid matching for recognizing natural scene categories. In *CVPR*, 2006.
- Wenbin Li. Learning multi-scale representations for material classification. In Xiaoyi Jiang, Joachim Hornegger, and Reinhard Koch, editors, *Pattern Recognition*, volume 8753, pages 757–764. Springer International Publishing, 2014.
- Julien Mairal, Piotr Koniusz, Zaid Harchaoui, and Cordelia Schmid. Convolutional kernel networks. In *NIPS*. 2014.
- Timo Ojala, Matti Pietikainen, and Topi Maenpaa. Multiresolution gray-scale and rotation invariant texture classification with local binary patterns. *TPAMI*, 24(7), 2002.
- Xiaojiang Peng, Limin Wang, Yu Qiao, and Qiang Peng. Boosting vlad with supervised dictionary learning and high-order statistics. In *ECCV*. 2014.
- Xavier Pennec, Pierre Fillard, and Nicholas Ayache. A riemannian framework for tensor computing. *IJCV*, 66(1), 2006.
- F. Perronnin and C. Dance. Fisher kernels on visual vocabularies for image categorization. In *CVPR*, 2007.
- Florent Perronnin, Jorge Sánchez, and Yan Liu. Large-scale image categorization with explicit data embedding. In *CVPR*, 2010a.
- Florent Perronnin, Jorge Sánchez, and Thomas Mensink. Improving the fisher kernel for large-scale image classification. In *ECCV*. 2010b.
- Ali Rahimi and Benjamin Recht. Random features for large-scale kernel machines. In *NIPS*, 2007.
- Walter Rudin. *Fourier analysis on groups*. John Wiley & Sons, 2011.
- Olga Russakovsky, Yuanqing Lin, Kai Yu, and Li Fei-Fei. Object-centric spatial pooling for image classification. In *ECCV*, pages 1–15. Springer, 2012.
- Lavanya Sharan, Ce Liu, Ruth Rosenholtz, and Edward H Adelson. Recognizing materials using perceptually inspired features. *IJCV*, 103(3), 2013.
- Terence Sim, Simon Baker, and Maan Bsat. The cmu pose, illumination, and expression database. *TPAMI*, 25(12), 2003.
- Josef Sivic, Bryan C Russell, Alexei A Efros, Andrew Zisserman, and William T Freeman. Discovering objects and their location in images. In *ICCV*, 2005.
- Suvrit Sra. A new metric on the manifold of kernel matrices with application to matrix geometric means. In *NIPS*, pages 144–152, 2012.

- Diego Tosato, Mauro Spera, Marco Cristani, and Vittorio Murino. Characterizing humans on riemannian manifolds. *TPAMI*, 35(8), 2013.
- Oncel Tuzel, Fatih Porikli, and Peter Meer. Region covariance: A fast descriptor for detection and classification. In *ECCV*. 2006.
- Oncel Tuzel, Fatih Porikli, and Peter Meer. Pedestrian detection via classification on Riemannian manifolds. *TPAMI*, 30(10), 2008.
- Jan C van Gemert, Cor J Veenman, Arnold WM Smeulders, and J-M Geusebroek. Visual word ambiguity. *TPAMI*, 32(7), 2010.
- A. Vedaldi and B. Fulkerson. VLFeat: An open and portable library of computer vision algorithms, 2008.
- A. Vedaldi and A. Zisserman. Efficient additive kernels via explicit feature maps. *TPAMI*, 34(3), 2012a.
- Andrea Vedaldi and Andrew Zisserman. Efficient additive kernels via explicit feature maps. *TPAMI*, 34(3), 2012b.
- Ruiping Wang, Huimin Guo, L.S. Davis, and Qionghai Dai. Covariance discriminative learning: A natural and efficient approach to image set classification. In *CVPR*, pages 2496–2503, June 2012.
- Yang Wang and Greg Mori. Human action recognition by semilattent topic models. *TPAMI*, 31(10), 2009.
- John Winn, Antonio Criminisi, and Thomas Minka. Object categorization by learned universal visual dictionary. In *ICCV*, 2005.
- Jianchao Yang, Kai Yu, Yihong Gong, and Thomas Huang. Linear spatial pyramid matching using sparse coding for image classification. In *CVPR*, 2009.
- Shaohua Kevin Zhou and Rama Chellappa. From sample similarity to ensemble similarity: Probabilistic distance measures in reproducing kernel hilbert space. *TPAMI*, 28(6), 2006.